



Graphs in Machine Learning

Michal Valko

Inria Lille - Nord Europe, France

TA: Pierre Perrault

Partially based on material by: Rob Fergus, Nikhil Srivastava,
Yiannis Koutis, Joshua Batson, Daniel Spielman



Last Lecture

- ▶ Inductive and transductive semi-supervised learning
- ▶ Manifold regularization
- ▶ Theory of Laplacian-based manifold methods
- ▶ Transductive learning stability based bounds
- ▶ Online semi-supervised learning
- ▶ Online incremental k -centers
- ▶ Examples of applications of online SSL
- ▶ Analysis of online SSL
- ▶ SSL Learnability
- ▶ When does graph-based SSL provably help?

This Lecture

- ▶ Scaling harmonic functions to millions of samples

Previous Lab Session

- ▶ 14. 11. 2018 by Pierre Perrault
- ▶ Content
 - ▶ Semi-supervised learning
 - ▶ Graph quantization
 - ▶ Offline face recognizer
- ▶ AR: **record a video with faces**
- ▶ Install VM (in case you have not done it yet for TD1)
- ▶ Short written report
- ▶ Questions to piazza
- ▶ **Deadline: 27. 11. 2017**

Next Lab Session/Lecture

- ▶ **DL for TD2: today**
- ▶ **No class or lab (TD) next week**
- ▶ 12.12.2018 by Pierre Perrault
- ▶ Content: Online and scalable algorithms
 - ▶ Online face recognizer
 - ▶ Iterative label propagation
 - ▶ Online k -centers
- ▶ AR: **record a video with faces**
- ▶ Short written report
- ▶ Questions to piazza
- ▶ **Deadline: 26.12.2018**

Final Class projects

- ▶ detailed description on the class website
- ▶ preferred option: you come up with the topic
- ▶ theory/implementation/review or a combination
- ▶ one or two people per project (exceptionally three)
- ▶ grade 60%: report + short presentation of the **team**
- ▶ deadlines
 - ▶ **21. 11. 2018** - strongly recommended DL for taking projects
 - ▶ 28. 11. 2018 - hard DL for taking projects
 - ▶ 07. 01. 2019 - submission of the project report
 - ▶ 11. 01. 2019 or later - project presentation
- ▶ list of suggested topics on piazza

Huge \mathcal{G}

when \mathcal{G} does not fit to memory

...or when we can't invert L

Scaling SSL with Graphs to Millions

Semi-supervised learning with graphs

$$\mathbf{f}^* = \min_{\mathbf{f} \in \mathbb{R}^N} (\mathbf{f} - \mathbf{y})^T \mathbf{C} (\mathbf{f} - \mathbf{y}) + \mathbf{f}^T \mathbf{L} \mathbf{f}$$

Let us see the same in eigenbasis of $\mathbf{L} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$, i.e., $\mathbf{f} = \mathbf{U}\boldsymbol{\alpha}$

$$\boldsymbol{\alpha}^* = \min_{\boldsymbol{\alpha} \in \mathbb{R}^N} (\mathbf{U}\boldsymbol{\alpha} - \mathbf{y})^T \mathbf{C} (\mathbf{U}\boldsymbol{\alpha} - \mathbf{y}) + \boldsymbol{\alpha}^T \boldsymbol{\Lambda} \boldsymbol{\alpha}$$

What is the problem with scalability?

Diagonalization of $N \times N$ matrix

What can we do? Let's take only first k eigenvectors $\mathbf{f} = \mathbf{U}\boldsymbol{\alpha}$!

Scaling SSL with Graphs to Millions

\mathbf{U} is now a $n \times k$ matrix

$$\alpha^* = \min_{\alpha \in \mathbb{R}^N} (\mathbf{U}\alpha - \mathbf{y})^T \mathbf{C}(\mathbf{U}\alpha - \mathbf{y}) + \alpha^T \mathbf{\Lambda} \alpha$$

Closed form solution is $(\mathbf{\Lambda} + \mathbf{U}^T \mathbf{C} \mathbf{U}) \alpha = \mathbf{U}^T \mathbf{C} \mathbf{y}$

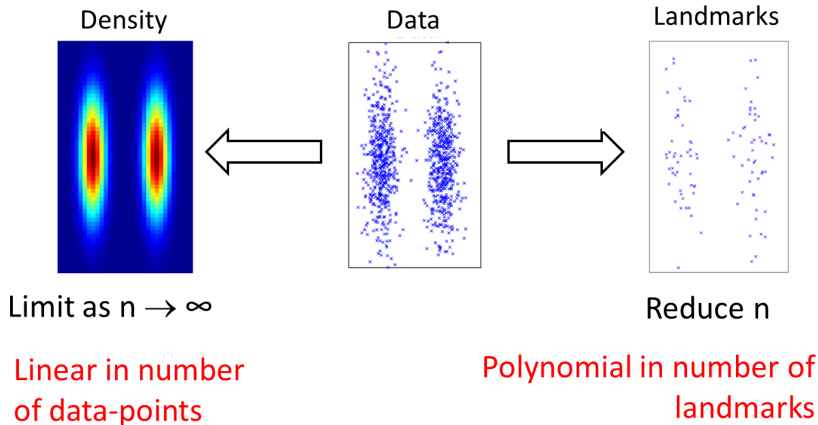
What is the size of this system of equation now?

Cool! Any problem with this approach?

Are there any reasonable assumptions when this is feasible?

Let's see what happens when $N \rightarrow \infty$!

Scaling SSL with Graphs to Millions



https://cs.nyu.edu/~fergus/papers/fwt_ssl.pdf

Scaling SSL with Graphs to Millions

What happens to \mathbf{L} when $N \rightarrow \infty$?

We have data $\mathbf{x}_i \in \mathbb{R}$ sampled from $p(\mathbf{x})$.

When $n \rightarrow \infty$, instead of vectors \mathbf{f} , we consider functions $F(\mathbf{x})$.

Instead of \mathbf{L} , we define \mathcal{L}_p - **weighted smoothness operator**

$$\mathcal{L}_p(F) = \frac{1}{2} \int (F(\mathbf{x}_1) - F(\mathbf{x}_2))^2 W(\mathbf{x}_1, \mathbf{x}_2) p(\mathbf{x}_1) p(\mathbf{x}_2) d\mathbf{x}_1 d\mathbf{x}_2$$

$$\text{with } W(\mathbf{x}_1, \mathbf{x}_2) = \frac{\exp(-\|\mathbf{x}_1 - \mathbf{x}_2\|^2)}{2\sigma^2}$$

\mathbf{L} defined the eigenvectors of increasing smoothness.

What defines \mathcal{L}_p ? **Eigenfunctions!**

Scaling SSL with Graphs to Millions

$$\mathcal{L}_p(F) = \frac{1}{2} \int (F(\mathbf{x}_1) - F(\mathbf{x}_2))^2 W(\mathbf{x}_1, \mathbf{x}_2) p(\mathbf{x}_1) p(\mathbf{x}_2) d\mathbf{x}_1 d\mathbf{x}_2$$

First eigenfunction

$$\Phi_1 = \underset{F: \int F^2(\mathbf{x}) p(\mathbf{x}) D(\mathbf{x}) d\mathbf{x} = 1}{\text{arg min}} \mathcal{L}_p(F)$$

where $D(\mathbf{x}) = \int_{\mathbf{x}_2} W(\mathbf{x}, \mathbf{x}_2) p(\mathbf{x}_2) d\mathbf{x}_2$

What is the solution? $\Phi_1(\mathbf{x}) = 1$ because $\mathcal{L}_p(1) = 0$

How to define Φ_2 ? same, constraining to be orthogonal to Φ_1

$$\int F(\mathbf{x}) \Phi_1(\mathbf{x}) p(\mathbf{x}) D(\mathbf{x}) d\mathbf{x} = 0$$

Scaling SSL with Graphs to Millions

Eigenfunctions of \mathcal{L}_p

Φ_3 as before, orthogonal to Φ_1 and Φ_2 etc.

How to define eigenvalues? $\lambda_k = \mathcal{L}_p(\Phi_k)$

Relationship to the discrete Laplacian

$$\frac{1}{N^2} \mathbf{f}^\top \mathbf{L} \mathbf{f} = \frac{1}{2N^2} \sum_{ij} W_{ij} (f_i - f_j)^2 \xrightarrow{N \rightarrow \infty} \mathcal{L}_p(F)$$

<http://www.informatik.uni-hamburg.de/ML/contents/people/luxburg/publications/>

Luxburg04_diss.pdf

<http://arxiv.org/pdf/1510.08110v1.pdf>

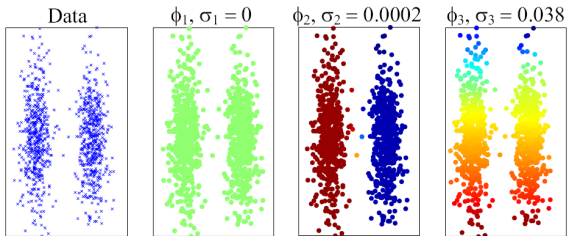
Isn't estimating eigenfunctions $p(\mathbf{x})$ more difficult?

Are there some "easy" distributions?

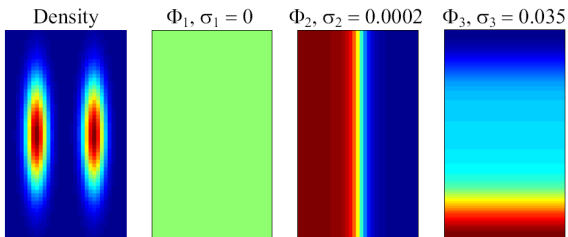
Can we compute it numerically?

Scaling SSL with Graphs to Millions

Eigenvectors



Eigenfunctions



Scaling SSL with Graphs to Millions

Factorized data distribution What if

$$p(\mathbf{s}) = p(s_1)p(s_2)\dots p(s_d)$$

In general, this is not true. But we can rotate data with $\mathbf{s} = \mathbf{R}\mathbf{x}$.



Treating each factor individually

$p_k \stackrel{\text{def}}{=} \text{marginal distribution of } s_k$

$\Phi_i(s_k) \stackrel{\text{def}}{=} \text{eigenfunction of } \mathcal{L}_{p_k} \text{ with eigenvalue } \lambda_i$

Then: $\Phi_i(\mathbf{s}) = \Phi_i(s_k)$ is eigenfunction of \mathcal{L}_p with λ_i

We only considered single-coordinate eigenfunctions.

Scaling SSL with Graphs to Millions

How to approximate 1D density? Histograms!

Algorithm of Fergus et al. [FWT09] for eigenfunctions

- ▶ Find \mathbf{R} such that $\mathbf{s} = \mathbf{R}\mathbf{x}$
- ▶ For each “independent” s_k approximate $p(s_k)$
- ▶ Given $p(s_k)$ numerically solve for eigensystem of \mathcal{L}_{p_k}

$$\left(\tilde{\mathbf{D}} - \widetilde{\mathbf{P}\mathbf{W}\mathbf{P}}\right) \mathbf{g} = \lambda \widehat{\mathbf{P}\mathbf{D}} \mathbf{g} \quad (\text{generalized eigensystem})$$

\mathbf{g} - vector of length $B \equiv$ number of bins

\mathbf{P} - density at discrete points

$\tilde{\mathbf{D}}$ - diagonal sum of the columns of $\widetilde{\mathbf{P}\mathbf{W}\mathbf{P}}$

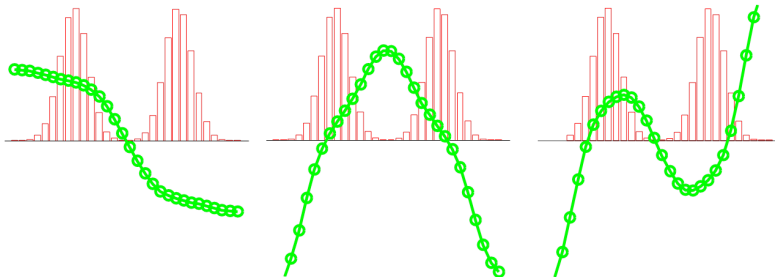
$\widehat{\mathbf{D}}$ - diagonal sum of the columns of $\widetilde{\mathbf{P}\mathbf{W}}$

- ▶ Order eigenfunctions by increasing eigenvalues

https://cs.nyu.edu/~fergus/papers/fwt_ssl.pdf

Scaling SSL with Graphs to Millions

Numerical 1D Eigenfunctions



1st Eigenfunction
of $h(x_1)$

2nd Eigenfunction
of $h(x_1)$

3rd Eigenfunction
of $h(x_1)$

https://cs.nyu.edu/~fergus/papers/fwt_ssl.pdf

Scaling SSL with Graphs to Millions

Computational complexity for $N \times d$ dataset

Typical harmonic approach

one diagonalization of $N \times N$ system

Numerical eigenfunctions with B bins and k eigenvectors

d eigenvector problems of $B \times B$

$$\left(\tilde{\mathbf{D}} - \mathbf{P}\tilde{\mathbf{W}}\mathbf{P}\right)\mathbf{g} = \lambda\mathbf{P}\hat{\mathbf{D}}\mathbf{g}$$

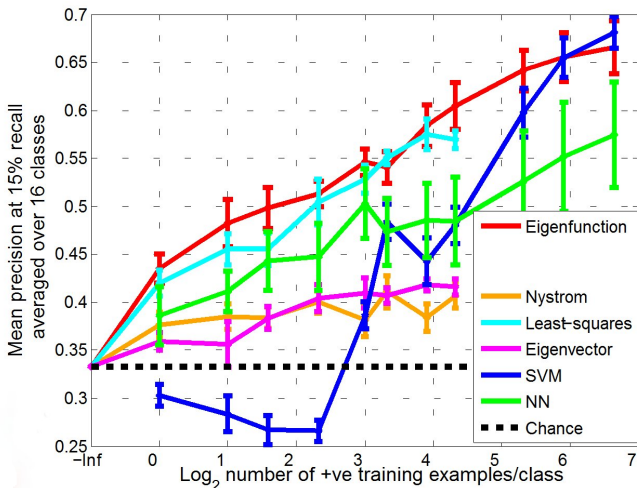
one $k \times k$ least squares problem

$$\left(\mathbf{\Lambda} + \mathbf{U}^T\mathbf{C}\mathbf{U}\right)\alpha = \mathbf{U}^T\mathbf{C}\mathbf{y}$$

some details: several approximation, eigenvectors only linear combinations single-coordinate eigenvectors, ...

When d is not too big then N can be in millions!

Scaling SSL with Graphs to Millions



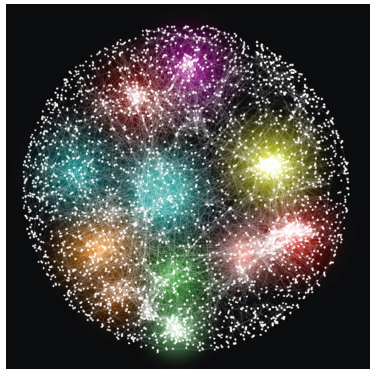
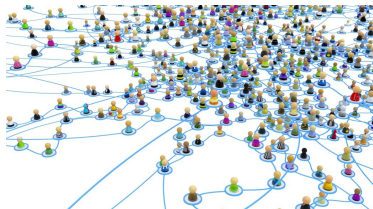
CIFAR experiments https://cs.nyu.edu/~fergus/papers/fwt_ssl.pdf

Sparsify \mathcal{G}

with no assumptions

...and we need to process is anyway

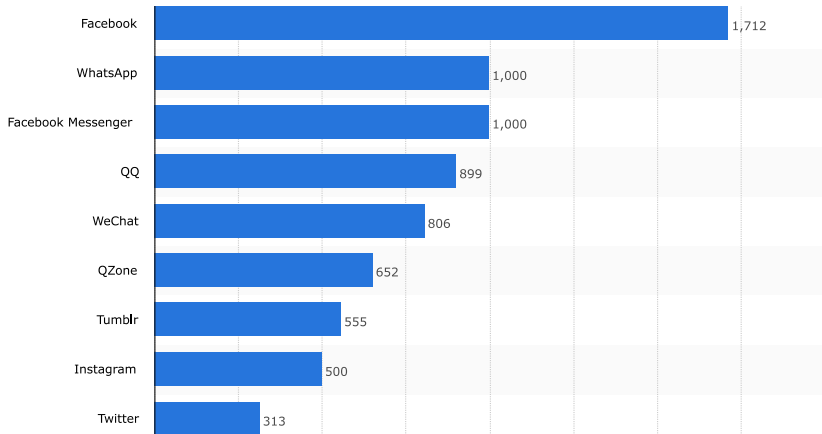
Large scale Machine Learning on Graphs



<http://blog.carsten-eickhoff.com>

Botstein et al.

Are we large yet?



"One **trillion** edges: graph processing at Facebook-scale."
Ching et al., VLDB 2015

Computational bottlenecks

In theory:

Space

$[\mathcal{O}(m), \mathcal{O}(n^2)]$ to store

Time

$\mathcal{O}(n^2)$ to construct
 $\mathcal{O}(n^3)$ to run algorithms

In practice:

- ▶ 2012 Common Crawl Corpus:
 - 3.5 Billion pages (45 GB)
 - 128 Billion edges (331 GB)
- ▶ Pagerank on Facebook Graph:
 - 3 minutes per iteration, hundreds of iterations, tens of hours on 200 machines, run once per day

Two phases

1 Preprocessing:

From vectorial data: Collect a dataset $\mathbf{X} \in \mathbb{R}^{n \times d}$, construct a graph \mathbf{G} using a similarity function

Prepare the graph: Need to check if graph is connected, make it directed/undirected, build Laplacian

Load it on the machine: On a single machine if possible, if not find smart way to distribute it

2 Run your algorithm on the graph

Large scale graph construction

Main bottleneck: **time**

- ▶ Constructing k -nn graph takes $\mathcal{O}(n^2 \log(n))$, too slow
- ▶ Constructing ε graph takes $\mathcal{O}(n^2)$, still too slow
- ▶ In both cases bottleneck is the same, given a node finding close nodes (k neighbours or ε neighbourhood)

Fundamental limit: just looking at all similarities already too slow.

Can we find close neighbours without checking all distances?

Distance Approximation

Split your data in small subset of close points

Can find efficiently some (not all) of the neighbours.

- ▶ Iterative Quantization
- ▶ KD-Trees – Cover Trees – NN search is $\mathcal{O}(\log N)$ per node
- ▶ Locality Sensitive Hashing (LSH)

More general problem: learning good codeword representation

Storing graph in memory

Main bottleneck: **space**.

As a Fermi (back-of-the-envelope) problem

- ▶ Storing a graph with m edges require to store m tuples $(i, j, w_{i,j})$ of 64 bit (8 bytes) doubles or int.
- ▶ For standard cloud providers, the largest compute-optimized instances has 36 cores, but only 60 GB of memory.
- ▶ We can store $60 * 1024^3 / (3 * 8) \sim 2.6 \times 10^9$ (2.6 billion) edges in a single machine memory.

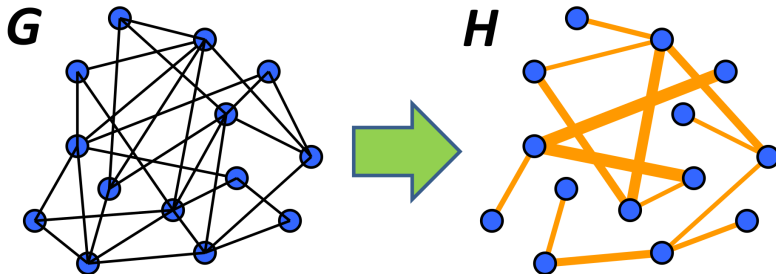
Storing graph in memory

But wait a minute

- ▶ Natural graphs are sparse.
 - ↳ For some it is true, for some it is false (e.g. Facebook average user has 300 friends, Twitter averages 208 followers)
Subcomponents are very dense, and they grow denser over time
- ▶ I will construct my graph sparse
 - ↳ Losing large scale relationship, losing regularization
- ▶ I will split my graph across multiple machines
 - ↳ Your algorithm does not know that.
What if it needs nonlocal data? Iterative algorithms?
More on this later

Graph Sparsification

Goal: Get graph G and find sparse H



Graph Sparsification: What is sparse?

What does **sparse** graph mean?

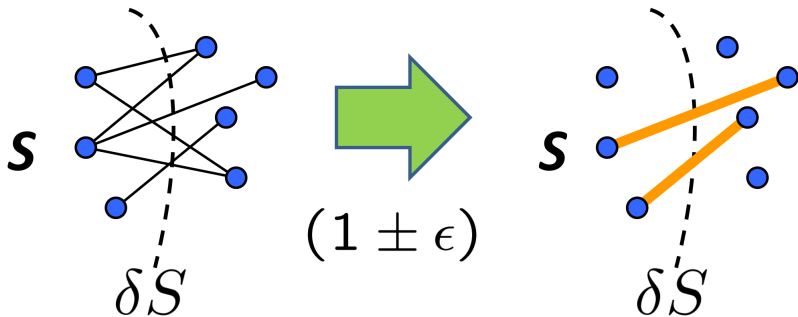
- ▶ average degree < 10 is pretty sparse
- ▶ for billion nodes even 100 should be ok
- ▶ in general: average degree $< \text{polylog } n$

Are all edges important?

in a tree — sure, in a dense graph perhaps not

Graph Sparsification: What is **good** sparse?

Good sparse by Benczúr and Karger (1996) = **cut preserving!**



H approximates G well iff $\forall S \subset V$, sum of edges on δS remains

δS = edges leaving S

<https://math.berkeley.edu/~nikhil/>

Graph Sparsification: What is **good** sparse?

Good sparse by Benczúr and Karger (1996) = **cut preserving!**

Why did they care? faster mincut/maxflow

Recall what is a cut: $\text{cut}_G(S) = \sum_{i \in S, j \in \bar{S}} w_{i,j}$

Define G and H are $(1 \pm \varepsilon)$ -**cut similar** when $\forall S$

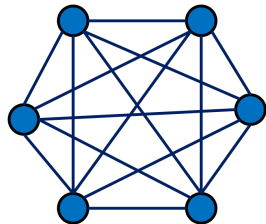
$$(1 - \varepsilon)\text{cut}_H(S) \leq \text{cut}_G(S) \leq (1 + \varepsilon)\text{cut}_H(S)$$

Is this always possible? Benczúr and Karger (1996): Yes!

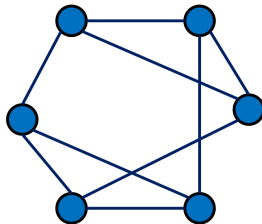
$\forall \varepsilon \exists (1 + \varepsilon)$ -cut similar \tilde{G} with $\mathcal{O}(n \log n / \varepsilon^2)$ edges s.t. $E_H \subseteq E$
and computable in $\mathcal{O}(m \log^3 n + m \log n / \varepsilon^2)$ time n nodes, m edges

Graph Sparsification: What is **good** sparse?

$G = K_n$



$H = d$ -regular (random)



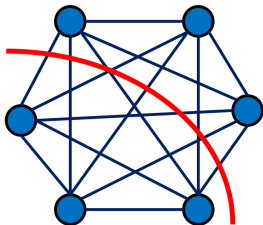
How many edges?

$$|E_G| = \mathcal{O}(n^2)$$

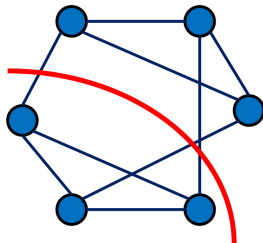
$$|E_H| = \mathcal{O}(dn)$$

Graph Sparsification: What is **good** sparse?

$G = K_n$



$H = d$ -regular (random)



What are the cut weights for any S ?

$$w_G(\delta S) = |S| \cdot |\bar{S}|$$

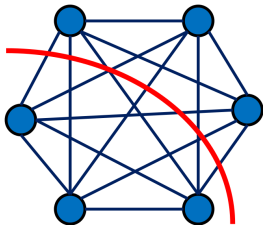
$$w_H(\delta S) \approx \frac{d}{n} \cdot |S| \cdot |\bar{S}|$$

$$\forall S \subset V : \frac{w_G(\delta S)}{w_H(\delta S)} \approx \frac{n}{d}$$

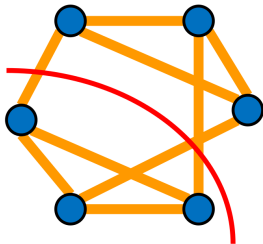
Could be large :(What to do?

Graph Sparsification: What is **good** sparse?

$G = K_n$



$H = d$ -regular (random)



What are the cut weights for any S ?

$$w_G(\delta S) = |S| \cdot |\bar{S}| \qquad w_H(\delta S) \approx \frac{d}{n} \cdot \frac{n}{d} \cdot |S| \cdot |\bar{S}|$$

$$\forall S \subset V : \frac{w_G(\delta S)}{w_H(\delta S)} \approx 1$$

Benczúr & Karger: Can find such H quickly for any G !

Graph Sparsification: What is **good** sparse?

Recall if $\mathbf{f} \in \{0, 1\}^n$ represents S then $\mathbf{f}^\top \mathbf{L}_G \mathbf{f} = \text{cut}_G(S)$

$$(1 - \varepsilon) \text{cut}_H(S) \leq \text{cut}_G(S) \leq (1 + \varepsilon) \text{cut}_H(S)$$

becomes

$$(1 - \varepsilon) \mathbf{f}^\top \mathbf{L}_H \mathbf{f} \leq \mathbf{f}^\top \mathbf{L}_G \mathbf{f} \leq (1 + \varepsilon) \mathbf{f}^\top \mathbf{L}_H \mathbf{f}$$

If we ask this only for $\mathbf{f} \in \{0, 1\}^n \rightarrow (1 + \varepsilon)$ -cut similar combinatorial
Benczúr & Karger (1996)

If we ask this for all $\mathbf{f} \in \mathbb{R}^n \rightarrow (1 + \varepsilon)$ -spectrally similar
Spielman & Teng (2004)

Spectral sparsifiers are stronger!

but checking for spectral similarity is easier

Spectral Graph Sparsification

Rayleigh-Ritz gives:

$$\lambda_{\min} = \min \frac{\mathbf{x}^T \mathbf{L} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \quad \text{and} \quad \lambda_{\max} = \max \frac{\mathbf{x}^T \mathbf{L} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$$

What can we say about $\lambda_i(G)$ and $\lambda_i(H)$?

$$(1 - \varepsilon) \mathbf{f}^T \mathbf{L}_G \mathbf{f} \leq \mathbf{f}^T \mathbf{L}_H \mathbf{f} \leq (1 + \varepsilon) \mathbf{f}^T \mathbf{L}_G \mathbf{f}$$

Eigenvalues are approximated well!

$$(1 - \varepsilon) \lambda_i(G) \leq \lambda_i(H) \leq (1 + \varepsilon) \lambda_i(G)$$

Using matrix ordering notation $(1 - \varepsilon) \mathbf{L}_G \preceq \mathbf{L}_H \preceq (1 + \varepsilon) \mathbf{L}_G$

As a consequence, $\arg \min_{\mathbf{x}} \|\mathbf{L}_H \mathbf{x} - \mathbf{b}\| \approx \arg \min_{\mathbf{x}} \|\mathbf{L}_G \mathbf{x} - \mathbf{b}\|$

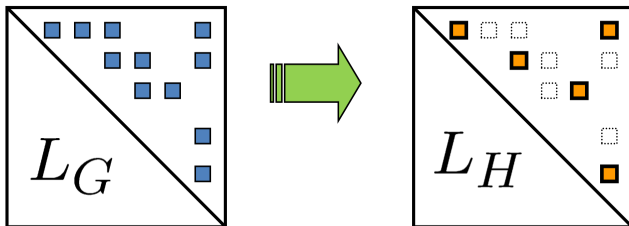
Spectral Graph Sparsification

Let us consider unweighted graphs: $w_{ij} \in \{0, 1\}$

$$\mathbf{L}_G = \sum_{ij} w_{ij} \mathbf{L}_{ij} = \sum_{ij \in E} \mathbf{L}_{ij} = \sum_{ij \in E} (\delta_i - \delta_j)(\delta_i - \delta_j)^T = \sum_{e \in E} \mathbf{b}_e \mathbf{b}_e^T$$

We look for a **subgraph** H

$$\mathbf{L}_H = \sum_{e \in E} s_e \mathbf{b}_e \mathbf{b}_e^T \quad \text{where } s_e \text{ is a new weight of edge } e$$



Spectral Graph Sparsification

We want $(1 - \varepsilon)\mathbf{L}_G \preceq \mathbf{L}_H \preceq (1 + \varepsilon)\mathbf{L}_G$

Equivalent, given $\mathbf{L}_G = \sum_{e \in E} \mathbf{b}_e \mathbf{b}_e^T$ find \mathbf{s} , s.t. $\mathbf{L}_G \preceq \sum_{e \in E} s_e \mathbf{b}_e \mathbf{b}_e^T \preceq \kappa \cdot \mathbf{L}_G$

Forget \mathbf{L} , given $\mathbf{A} = \sum_{e \in E} \mathbf{a}_e \mathbf{a}_e^T$ find \mathbf{s} , s.t. $\mathbf{A} \preceq \sum_{e \in E} s_e \mathbf{a}_e \mathbf{a}_e^T \preceq \kappa \cdot \mathbf{A}$

Same as, given $\mathbf{I} = \sum_{e \in E} \mathbf{v}_e \mathbf{v}_e^T$ find \mathbf{s} , s.t. $\mathbf{I} \preceq \sum_{e \in E} s_e \mathbf{v}_e \mathbf{v}_e^T \preceq \kappa \cdot \mathbf{I}$

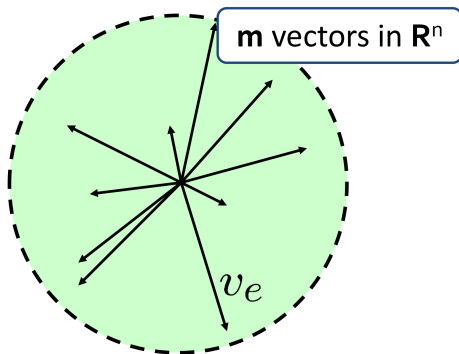
How to get it? $\mathbf{v}_e \leftarrow \mathbf{A}^{-1/2} \mathbf{a}_e$

Then $\sum_{e \in E} s_e \mathbf{v}_e \mathbf{v}_e^T \approx \mathbf{I} \iff \sum_{e \in E} s_e \mathbf{a}_e \mathbf{a}_e^T \approx \mathbf{A}$

multiplying by $\mathbf{A}^{1/2}$ on both sides

Spectral Graph Sparsification: Intuition

How does $\sum_{e \in E} \mathbf{v}_e \mathbf{v}_e^T = \mathbf{I}$ look like geometrically?



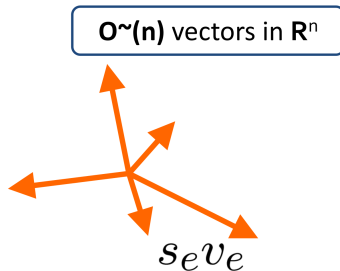
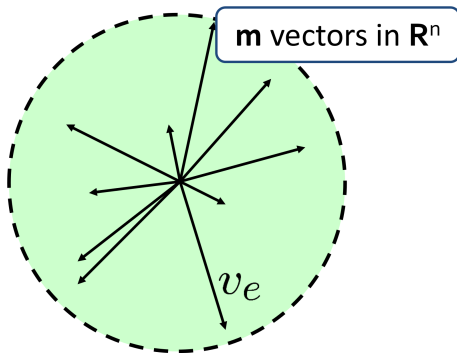
Decomposition of identity: $\forall \mathbf{u}$ (unit vector): $\sum_{e \in E} (\mathbf{u}^T \mathbf{v}_e)^2 = 1$

moment ellipse is a sphere

<https://math.berkeley.edu/~nikhil/>

Spectral Graph Sparsification: Intuition

What are we doing by choosing H ?

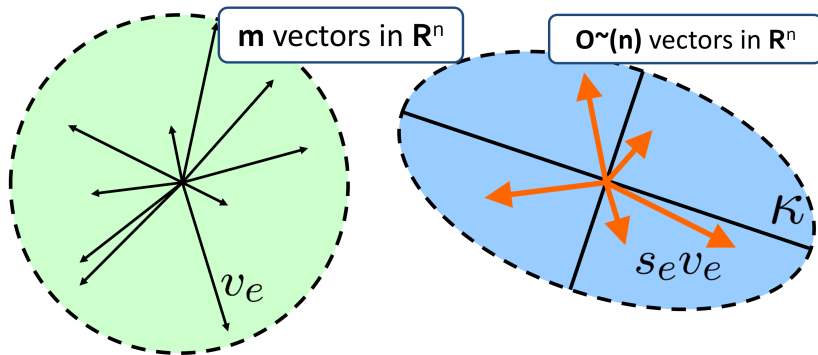


We take a subset of these e_e s and scale them!

<https://math.berkeley.edu/~nikhil/>

Spectral Graph Sparsification: Intuition

What kind of scaling do we want?



Such that the blue ellipsoid looks like identity!

the blue eigenvalues are between 1 and κ

<https://math.berkeley.edu/~nikhil/>

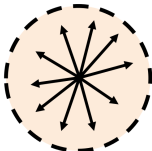
Spectral Graph Sparsification: Intuition

Example: What happens with K_n ?

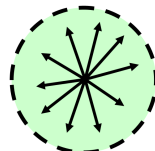
K_n graph



$$\sum_{e \in E} \mathbf{b}_e \mathbf{b}_e^T = \mathbf{L}_G$$



$$\sum_{e \in E} \mathbf{v}_e \mathbf{v}_e^T = \mathbf{I}$$



It is already isotropic! (looks like a sphere)

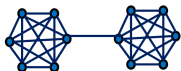
rescaling $\mathbf{v}_e = \mathbf{L}^{-1/2} \mathbf{b}_e$ does not change the shape

<https://math.berkeley.edu/~nikhil/>

Spectral Graph Sparsification: Intuition

Example: What happens with a dumbbell?

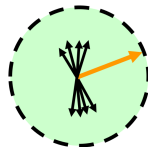
K_n graph



$$\sum_{e \in E} \mathbf{b}_e \mathbf{b}_e^T = \mathbf{L}_G$$



$$\sum_{e \in E} \mathbf{v}_e \mathbf{v}_e^T = \mathbf{I}$$



The vector corresponding to the link gets stretched!

because this transformation makes all the directions important

rescaling reveals the vectors that are critical

<https://math.berkeley.edu/~nikhil/>

Spectral Graph Sparsification: Intuition

What is this rescaling $\mathbf{v}_e = \mathbf{L}_G^{-1/2} \mathbf{b}_e$ doing to the norm?

$$\|\mathbf{v}_e\|^2 = \left\| \mathbf{L}_G^{-1/2} \mathbf{b}_e \right\|^2 = \mathbf{b}_e^\top \mathbf{L}_G^{-1} \mathbf{b}_e = R_{\text{eff}}(e)$$

reminder $R_{\text{eff}}(e)$ is the potential difference between the nodes when injecting a unit current

In other words: $R_{\text{eff}}(e)$ is related to the edge importance!

Electrical intuition: We want to find an electrically similar H and the importance of the edge is its effective resistance $R_{\text{eff}}(e)$.

Edges with higher R_{eff} are more **electrically significant!**

Spectral Graph Sparsification

Todo: Given $\mathbf{I} = \sum_e \mathbf{v}_e \mathbf{v}_e^T$, find a sparse reweighting.

Randomized algorithm that finds \mathbf{s} :

- ▶ Sample $n \log n / \varepsilon^2$ with replacement $p_i \propto \|\mathbf{v}_e\|^2$ (resistances)
- ▶ Reweigh: $s_i = 1/p_i$ (to be unbiased)

Does this work?

Application of Matrix Chernoff Bound by Rudelson (1999)

$$1 - \varepsilon \prec \lambda \left(\sum_e s_e \mathbf{v}_e \mathbf{v}_e^T \right) \prec 1 + \varepsilon$$

finer bounds now available

What is the the biggest problem here? Getting the p_i s!

Spectral Graph Sparsification

We want to make this algorithm fast.

How can we compute the effective resistances?

Solve a linear system $\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{L}_G \mathbf{x} - \mathbf{b}_e\|$ and then $R_{\text{eff}} = \mathbf{b}_e^T \hat{\mathbf{x}}$

Gaussian Elimination $\mathcal{O}(n^3)$

Fast Matrix Multiplication $\mathcal{O}(n^{2.37})$

Spielman & Teng (2004) $\mathcal{O}(m \log^{30} n)$

Koutis, Miller, and Peng (2010) $\mathcal{O}(m \log n)$

► Fast solvers for SDD systems:

↳ use sparsification internally

all the way until you hit the turtles

still infeasible when m is large

Spectral Graph Sparsification

Chicken and egg problem

We need R_{eff} to compute a sparsifier H \leftarrow

\hookrightarrow We need a sparsifier H to compute R_{eff}

Sampling according to approximate effective resistances

$R_{\text{eff}} \leq \tilde{R}_{\text{eff}} \leq \alpha R_{\text{eff}}$ give approximate sparsifier $\mathbf{L}_G \preceq \mathbf{L}_H \preceq \alpha \kappa \mathbf{L}_G$

Start with very poor approximation \tilde{R}_{eff} and poor sparsifier.

Use \tilde{R}_{eff} to compute an improved approximate sparsifier H \leftarrow

\hookrightarrow Use the sparsifier H to compute improved approximate \tilde{R}_{eff}

Computing \tilde{R}_{eff} using the sparsifier is fast ($m = \mathcal{O}(n \log(n))$), and not too many iterations are necessary.

What can I use sparsifiers for?

- ▶ Graph linear systems: minimum cut, maximum flow, Laplacian regression, SSL
- ▶ More in general, solving Strongly Diagonally Dominant (SDD) linear systems
 - ↳ electric circuit, fluid equations, finite elements methods
- ▶ Various embeddings: k-means, spectral clustering.

But what if my problems have no use for spectral guarantees?

Or if my boss does not trust approximation methods

Distributed graph processing

Large graphs do not fit in memory

Get more memory

- ↳ Either slower but larger memory
Or fast memory but divided among many machines

Many challenges

Needs to be scalable

- ↳ minimize pass over data / communication costs

Needs to be consistent

- ↳ updates should propagate properly

Distributed graph processing

Different choices have different impacts: for example splitting the graph according to nodes or according to edges.

Many computation models (academic and commercial) each with its pros and cons

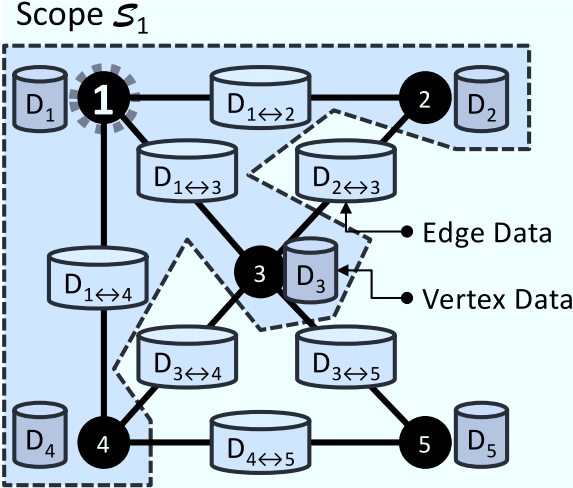
MapReduce

MPI

Pregel

Graphlab

The GraphLab abstraction



The GraphLab abstraction

```
In [1]: import sframe
```

```
In [2]: edges = sframe.SFrame.read_csv('/media/sf_share/td3_example_edges.csv')
```

```
In [3]: vertices = sframe.SFrame.read_csv('/media/sf_share/td3_example_vertices.csv')
```

```
In [4]: G = sframe.SGraph(edges= edges, vertices=vertices, src_field='src', dst_field='dst')
```

```
In [5]: G
```

```
Out[5]: SGraph({'num_edges': 26, 'num_vertices': 9})  
Vertex Fields:['__id', 'f']  
Edge Fields:['__src_id', '__dst_id', 'weight']
```

The GraphLab abstraction

Under the hood: tabular representation

Columns:
__id int
f float

Rows: 9

Data:

__id	f
5	0.51
7	0.82
10	0.08
2	0.82
6	0.85
9	0.83
3	0.18
1	0.35
4	0.36

[9 rows x 2 columns]

Columns:
__src_id int
__dst_id int
weight float

Rows: 26

Data:

__src_id	__dst_id	weight
7	5	0.13185
5	7	0.13185
7	7	0.026779
10	7	0.57121
7	10	0.57121
10	2	0.94047
7	6	0.64528
5	3	0.93374
10	3	0.31713
5	1	0.57796

[26 rows x 3 columns]

Note: Only the head of the SFrame is printed.

The GraphLab abstraction

```
In [1]: import sframe
```

```
In [2]: G = sframe.SGraph()
```

```
In [3]: G
```

```
Out[3]: SGraph({'num_edges': 0, 'num_vertices': 0})  
Vertex Fields:['__id']  
Edge Fields:['__src_id', '__dst_id']
```

```
In [1]: import sframe
```

```
In [2]: G = sframe.SGraph()
```

```
In [3]: G
```

```
Out[3]: SGraph({'num_edges': 0, 'num_vertices': 0})  
Vertex Fields:['__id']  
Edge Fields:['__src_id', '__dst_id']
```

```
In [4]: G.add_edges(sframe.Edge(1,2))
```

```
Out[4]: SGraph({'num_edges': 1, 'num_vertices': 2})  
Vertex Fields:['__id']  
Edge Fields:['__src_id', '__dst_id']
```

The GraphLab abstraction

- ▶ The graph is immutable. why?
- ▶ All computations are executed asynchronously
 - ↳ We do not know the order of execution
We do not even know where the node is stored
what data can we access?
- ▶ The data is stored in the graph itself
 - ↳ only access local data
- ▶ Functional programming approach

The GraphLab abstraction

```
triple_apply(triple_apply_fn, mutated_fields, input_fields=None)
```

processes all edges asynchronously and in parallel

```
>>> PARALLEL FOR (source, edge, target) AS triple in G:  
...     LOCK (triple.source, triple.target)  
...     (source, edge, target) = triple_apply_fn(triple)  
...     UNLOCK (triple.source, triple.target)  
... END PARALLEL FOR
```

- ▶ No guarantees on order of execution
- ▶ Updating (src,edge,dst) would violate immutability
- ▶ triple_apply_fn receives a copy of (src,edge,dst)
 - ↳ returns an updated (src',edge',dst')
 - use return values to build a new graph

The GraphLab abstraction

triple_apply_fn is a pure function

Function in the mathematical sense, same input gives same output.

```
1 def triple_apply_fn(src, edge, dst):
2     #can only access data stored in src, edge, and dst,
3     #three dictionaries containing a copy of the
4     #fields indicated in mutated_fields
5     f = dst['f']
6
7     #inputs are copies, this does not change original edge
8     edge['weight'] = g(f)
9
10    return ({'f': dst['f']}, edge, dst)
```

The GraphLab abstraction

An example, computing degree of nodes

```
1 def degree_count_fn (src , edge , dst):  
2     src ['degree'] += 1  
3     dst ['degree'] += 1  
4     return (src , edge , dst)  
5  
6 G_count = G.triple_apply (degree_count_fn , 'degree')
```

The GraphLab abstraction

Slightly more complicated example, suboptimal pagerank

```
1 #assume each node in G has a field 'degree' and 'pagerank'
2 #initialize 'pagerank' = 1/n for all nodes
3
4 def weight_count_fn (src , edge , dst):
5     dst['degree'] += edge['weight']
6     return (src , edge , dst)
7
8 def pagerank_step_fn (src , edge , dst):
9     dst['pagerank'] += (edge['weight'] * src['pagerank']
10                        /dst['degree'])
11     return (src , edge , dst)
12
13 G_pagerank = G.triple_apply(weight_count_fn , 'degree')
14
15 while not converged(G_pagerank):
16     G_pagerank = G_pagerank.triple_apply(
17         pagerank_step_fn , 'pagerank')
```

How many iterations to convergence?

Graph Spectral Sparsification

Definition ([SS11])

An ε -sparsifier of \mathcal{G} is a reweighted subgraph \mathcal{H} whose Laplacian $\mathbf{L}_{\mathcal{H}}$ satisfies

$$(1 - \varepsilon)\mathbf{L}_{\mathcal{G}} \preceq \mathbf{L}_{\mathcal{H}} \preceq (1 + \varepsilon)\mathbf{L}_{\mathcal{G}} \quad (1)$$

Proposition ([SS11; Kyn+16])

There exists an algorithm that can construct an ε -sparsifier

- ▶ with only $\mathcal{O}(n \log(n)/\varepsilon^2)$ edges
- ▶ in $\mathcal{O}(m \log^2(n))$ time and $\mathcal{O}(n \log(n)/\varepsilon^2)$ space
- ▶ a single pass over the data

Graph Spectral Sparsification in Machine Learning

Laplacian smoothing (denoising): given $\mathbf{y} \triangleq \mathbf{f}^* + \xi$ and \mathcal{G} compute

$$\min_{\mathbf{f} \in \mathbb{R}^n} (\mathbf{f} - \mathbf{y})^\top (\mathbf{f} - \mathbf{y}) + \lambda \mathbf{f}^\top \mathbf{L}_{\mathcal{G}} \mathbf{f} \quad (2)$$

	Preproc	Time	Space
$\hat{\mathbf{f}} = (\lambda \mathbf{L}_{\mathcal{G}} + \mathbf{I})^{-1} \mathbf{y}$	0	$\mathcal{O}(m \log(n))$	$\mathcal{O}(m)$
$\tilde{\mathbf{f}} = (\lambda \mathbf{L}_{\mathcal{H}} + \mathbf{I})^{-1} \mathbf{y}$	$\mathcal{O}(m \log^2(n))$	$\mathcal{O}(n \log^2(n))$	$\mathcal{O}(n \log(n))$

Large computational improvement

↳ accuracy guarantees! [SWT16]

Need to approximate spectrum only up to regularization level λ

Ridge Graph Spectral Sparsification

Definition

An (ε, γ) -sparsifier of \mathcal{G} is a reweighted subgraph \mathcal{H} whose Laplacian $\mathbf{L}_{\mathcal{H}}$ satisfies

$$(1 - \varepsilon)\mathbf{L}_{\mathcal{G}} - \varepsilon\gamma\mathbf{I} \preceq \mathbf{L}_{\mathcal{H}} \preceq (1 + \varepsilon)\mathbf{L}_{\mathcal{G}} + \varepsilon\gamma\mathbf{I} \quad (3)$$

Mixed multiplicative/additive error

- ▶ large (i.e. $\geq \gamma$) directions reconstructed accurately
- ▶ small (i.e. $\leq \gamma$) directions uniformly approximated ($\gamma\mathbf{I}$)

Adapted from Randomized Linear Algebra (RLA) community

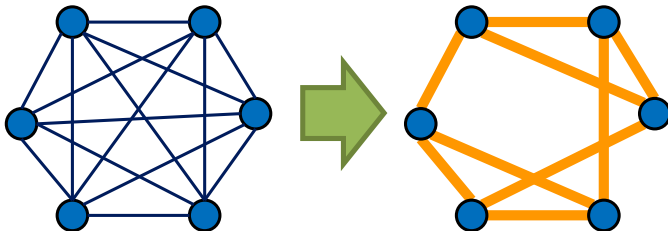
↳ PSD matrix low-rank approx. [AM15]

RLA → Graph: Improve over $\mathcal{O}(n \log n)$ exploiting regularization

Graph → RLA: Exploit $\mathbf{L}_{\mathcal{G}}$ structure for fast (ε, γ) -sparsification

How to construct an ε -sparsifier

For complete graphs, sample $\mathcal{O}(n \log(n))$ edges uniformly and reweight

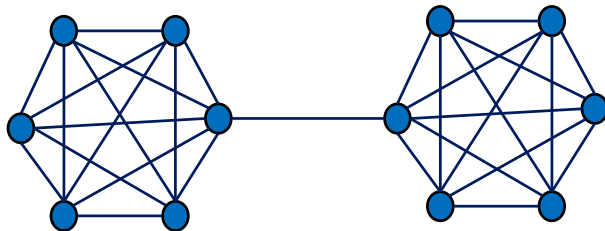


How to construct an ε -sparsifier

For generic graphs, sample $\mathcal{O}(n \log(n))$ edges uniformly?

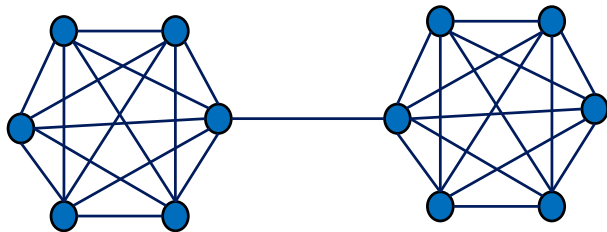
How to construct an ε -sparsifier

For generic graphs, sample $\mathcal{O}(n \log(n))$ edges uniformly?



How to construct an ε -sparsifier

For generic graphs, sample $\mathcal{O}(n \log(n))$ edges using effective resistance



Effective resistance $r_e = \mathbf{b}_e^T \mathbf{L}_G^+ \mathbf{b}_e$ of an edge

↳ inverse of number of alternative paths

↳ sum of r_e is $n - 1$

<https://math.berkeley.edu/~nikhil/>

How to construct an (ε, γ) -sparsifier

Definition

γ -effective resistance: $r_e(\gamma) = \mathbf{b}_e^\top (\mathbf{L}_G + \gamma \mathbf{I})^{-1} \mathbf{b}_e$

Effective dim.: $\mathbf{d}_{\text{eff}}(\gamma) = \sum_e r_e(\gamma) = \sum_{i=1}^n \frac{\lambda_i(\mathbf{L}_G)}{\lambda_i(\mathbf{L}_G) + \gamma} \leq n$

Can still be computed using fast graph solvers

↳ interpretation as inverse of alternative paths lost

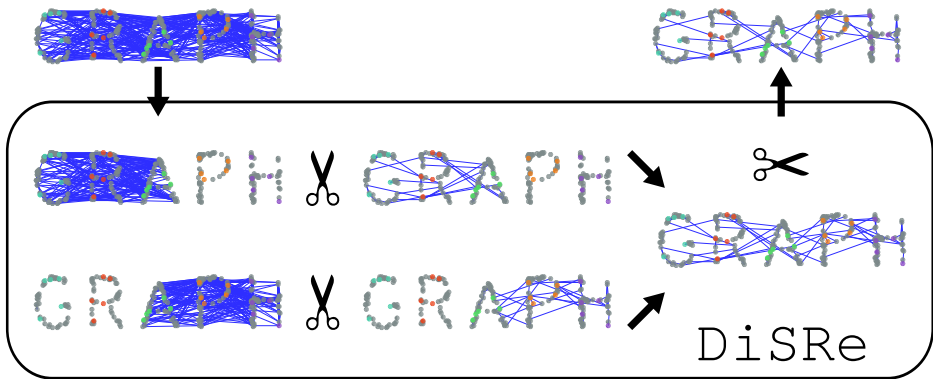
Most existing graph algorithms inapplicable [Kyn+16]

Most existing RLA algorithms too slow [CMM17]

Adapt SOA algorithm for kernel matrix approximation

SQUEAK, [CLV17]

DisRe



arbitrarily split in subgraphs that fit in a single machine
recursively merge-and-reduce until one graph left

↳ additive error cumulates!

↳ merge-and-resparsify

Sparsification



Compute $\tilde{p}_e^{(1)} \propto \tilde{r}_e^{(1)}(\gamma)$ using fast graph solver

For each edge e sample with probability $\tilde{p}_e^{(1)}$

w.h.p. (ϵ, γ) -accurate and use only

$\mathcal{O}(\text{diff}(\gamma) \log(n)) \leq \mathcal{O}(n \log(n))$ space

Sparsification



Compute $\tilde{p}_e^{(1)} \propto \tilde{r}_e^{(1)}(\gamma)$ using fast graph solver

For each edge e sample with probability $\tilde{p}_e^{(1)}$

w.h.p. (ϵ, γ) -accurate and use only

$\mathcal{O}(d_{\text{eff}}(\gamma) \log(n)) \leq \mathcal{O}(n \log(n))$ space

Sparsification



Compute $\tilde{p}_e^{(1)} \propto \tilde{r}_e^{(1)}(\gamma)$ using fast graph solver

For each edge e sample with probability $\tilde{p}_e^{(1)}$

w.h.p. (ϵ, γ) -accurate and use only

$\mathcal{O}(d_{\text{eff}}(\gamma) \log(n)) \leq \mathcal{O}(n \log(n))$ space

Sparsification



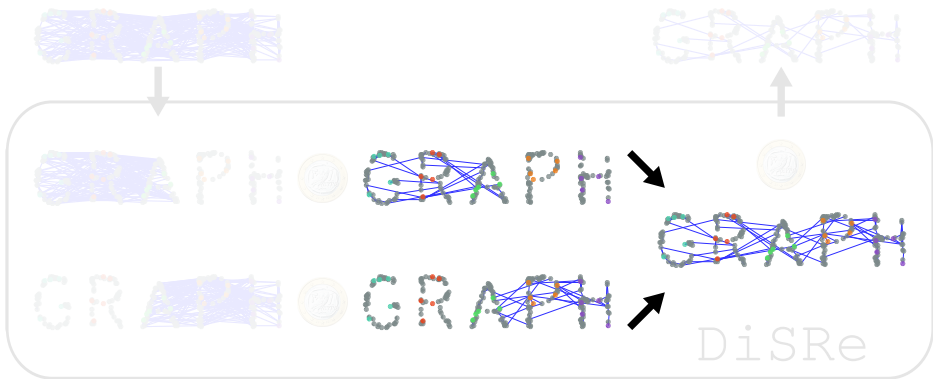
Compute $\tilde{p}_e^{(1)} \propto \tilde{r}_e^{(1)}(\gamma)$ using fast graph solver

For each edge e sample with probability $\tilde{p}_e^{(1)}$

w.h.p. (ϵ, γ) -accurate and use only

$\mathcal{O}(d_{\text{eff}}(\gamma) \log(n)) \leq \mathcal{O}(n \log(n))$ space

Merge



Combine sparsifiers, using $2\mathcal{O}(d_{\text{eff}}(\gamma) \log(n))$ space

twice as large as necessary

Merge-and-Resparsify

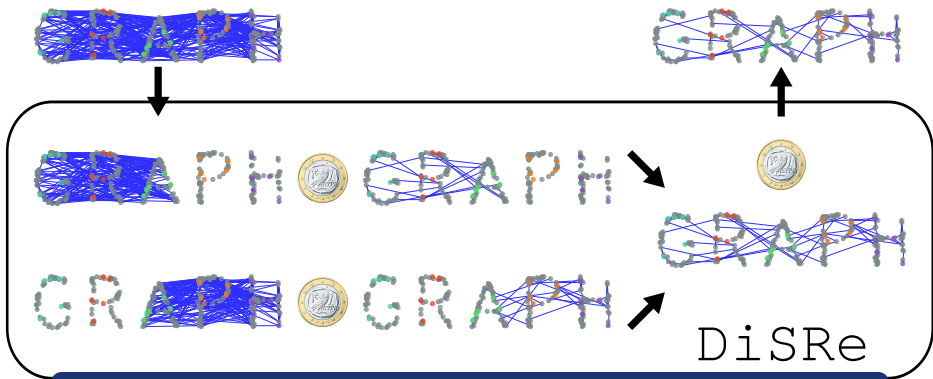


Compute $\tilde{p}_e^{(2)} \propto \min\{\tilde{r}_e^{(2)}(\gamma), \tilde{p}_e^{(1)}\}$ using fast graph solver

For each edge e sample with probability $\tilde{p}_e^{(2)} / \tilde{p}_e^{(1)}$

$$\text{survival probability } \frac{\tilde{p}_e^{(2)}}{\tilde{p}_e^{(1)}} \tilde{p}_e^{(1)}$$

DisRe guarantees

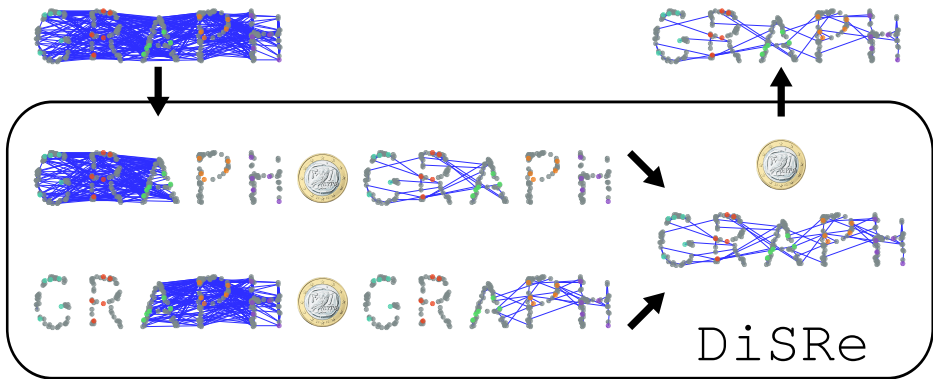


Theorem

Given an arbitrary graph \mathcal{G} w.h.p. DISRE satisfies

- (1) each sub-graphs is an (ϵ, γ) -sparsifier
- (2) with at most $\mathcal{O}(d_{\text{eff}}(\gamma) \log(n))$ edges.

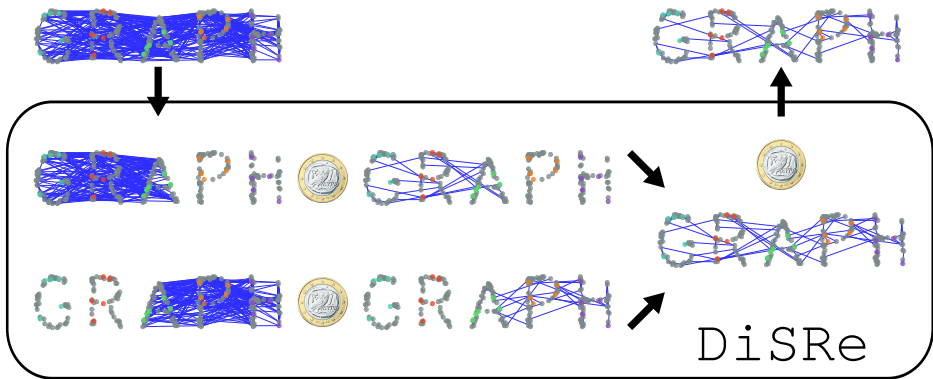
DisRe guarantees



Space: independent from m $\mathcal{O}(d_{\text{eff}}(\gamma) \log(n)) \leq \mathcal{O}(n \log(n))$

Time: $\mathcal{O}(d_{\text{eff}}(\gamma) \log^3(n))$ for fully balanced tree

DisRe guarantees



Communication: only $\mathcal{O}(\log(n))$ rounds

↳ removed edges are forgotten **single pass/streaming**

↳ point-to-point, centralization only to **choose tree**

Guarantees for Laplacian smoothing

$$\hat{\mathbf{f}} = (\lambda \mathbf{L}_G + \mathbf{I})^{-1} \mathbf{y}, \quad \tilde{\mathbf{f}} = (\lambda \mathbf{L}_H + \mathbf{I})^{-1} \mathbf{y}$$

Theorem ([SWT16] [CLV17])

If \mathbf{L}_H is an $(\varepsilon, 0)$ (ε, γ) -sparsifier of \mathbf{L}_G

$$\|\tilde{\mathbf{f}} - \hat{\mathbf{f}}\|_2^2 \leq \frac{\varepsilon^2}{1 - \varepsilon} (0.25 + \lambda\gamma) \left(\lambda \hat{\mathbf{f}}^\top \mathbf{L}_G \hat{\mathbf{f}} + \lambda\gamma \|\hat{\mathbf{f}}\|_2^2 \right).$$

$\mathcal{O}(d_{\text{eff}}(\gamma) \log(n))$ space, $\mathcal{O}(d_{\text{eff}}(\gamma) \log^3(n))$ time

↳ exploit regularization: \mathcal{H} sub-linear in n

Recover bound for ε -sparsifier when $\gamma \rightarrow 0$

↳ freely cross-validate γ since $d_{\text{eff}}(0) \leq n$

↳ trade-off between smoothness and decay of \mathbf{L}_G

Experiments

Dataset: Amazon co-purchase graph [YL15]

↳ **natural**, artificially sparse (true graph known only to Amazon)
↳ we compute 4-step random walk to recover removed co-purchases [GM15]

Target: eigenvector \mathbf{v} associated with $\lambda_2(\mathbf{L}_G)$ [SWT16]

$n = 334,863$ nodes, $m = 98,465,352$ edges (294 avg. degree)

Alg.	Parameters	$ \mathcal{E} $ ($\times 10^6$)	$\ \tilde{\mathbf{f}} - \mathbf{v}\ _2^2$ ($\sigma = 10^{-3}$)	$\ \tilde{\mathbf{f}} - \mathbf{v}\ _2^2$ ($\sigma = 10^{-2}$)
EX-ACT		98.5	0.067 ± 0.0004	0.756 ± 0.006
kN	$k = 60$	15.7	0.172 ± 0.0004	0.822 ± 0.002
DISRE	$\gamma = 0$	22.8	0.068 ± 0.0004	0.756 ± 0.005
DISRE	$\gamma = 10^2$	11.8	0.068 ± 0.0002	0.772 ± 0.004

Time: Loading \mathcal{G} from disk 90sec, DISRE 120sec ($k = 4 \times 32$ CPU), computing $\tilde{\mathbf{f}}$ 120sec, computing $\hat{\mathbf{f}}$ 720sec

Recap and open questions

Remark ([SWT16])

To the best of our knowledge, [graph sparsification] applications in machine learning have not yet been thoroughly pursued.

- ▶ introduction of (ε, γ) -sparsifiers to Graph ML
- ▶ DISRE, new distributed algorithm to construct (ε, γ) -sparsifiers
- ▶ new results for fast Laplacian Smoothing
- ▶ new results for fast SSL using ε -sparsifiers (at poster #76)

Open questions

- ▶ other accelerated Graph ML algorithms using (ε, γ) -sparsifiers
- ▶ more experiments on **dense** graphs
- ▶ Facebook: 300 average friends [Pew Research Center 2013]
- ▶ Twitter 453 average followers, 3.4x denser 2012-16 [LKF07]

Bibliography I



Ahmed El Alaoui and Michael W. Mahoney. “Fast randomized kernel methods with statistical guarantees”. In: *Neural Information Processing Systems*. 2015.



Daniele Calandriello, Alessandro Lazaric, and Michal Valko. “Distributed adaptive sampling for kernel matrix approximation”. In: *International Conference on Artificial Intelligence and Statistics*. 2017.



Michael B. Cohen, Cameron Musco, and Christopher Musco. “Input sparsity time low-rank approximation via ridge leverage score sampling”. In: *Symposium on Discrete Algorithms*. 2017.

Bibliography II



Rob Fergus, Yair Weiss, and Antonio Torralba. “Semi-Supervised Learning in Gigantic Image Collections”. In: *Neural Information Processing Systems*. 2009.



David F Gleich and Michael W Mahoney. “Using Local Spectral Methods to Robustify Graph-Based Learning Algorithms”. In: *Knowledge Discovery and Data Mining*. 2015.



Rasmus Kyng et al. “A Framework for Analyzing Resparsification Algorithms”. In: *Symposium on Discrete Algorithms*. 2016.

Bibliography III



Jure Leskovec, Jon Kleinberg, and Christos Faloutsos. “Graph evolution: Densification and shrinking diameters”. In: *Knowledge Discovery from Data* (Mar. 2007).



Daniel A Spielman and Nikhil Srivastava. “Graph sparsification by effective resistances”. In: *Journal on Computing* 40.6 (2011).



Veeranjaneyulu Sadhanala, Yu-xiang Wang, and Ryan J Tibshirani. “Graph Sparsification Approaches for Laplacian Smoothing”. In: *International Conference on Artificial Intelligence and Statistics*. 2016.

Bibliography IV



Jaewon Yang and Jure Leskovec. “Defining and evaluating network communities based on ground-truth”. In: *Knowledge and Information Systems* (2015).

Michal Valko

michal.valko@inria.fr

ENS Paris-Saclay, MVA 2018/2019

Sequel team, Inria Lille — Nord Europe

<https://team.inria.fr/sequel/>